



Reporte final

Eduardo Joel Cortez Valente A01746664

Ana Martínez Barbosa A01382889

José María Ibarra A01706970

Jorge Isidro Blanco Martínez A01745907

Maximiliano Benítez Ahumada A01752791

TC3006C. Inteligencia artificial avanzada para la ciencia de datos

15/09/2023

Contexto

La analítica de datos en la venta de bienes raíces puede proporcionar información valiosa tanto para los vendedores como para los compradores, ya que de esta manera se puede tener precios acordes a la propiedad, lo que hace que los compradores paguen lo justo tomando en cuenta las características de la casa y los vendedores pueden vender las propiedades en menor tiempo.

Las agencias inmobiliarias tendrán una clara ventaja competitiva frente a aquellas que no lo utilicen porque podrán hacer ventas más precisas proporcionando datos duros a los clientes, lo que los puede dejar más satisfechos y generar más ventas. Por lo tanto, La analítica de datos no solo aumenta la eficiencia y la precisión en la toma de decisiones, sino que también mejora la experiencia del cliente y proporciona información relevante sobre el mercado, lo cual puede ser crucial al momento de realizar una venta.

I. DESCRIPCIÓN BASE DE DATOS

Para la obtención de un modelo que permitiese generar predicciones acertadas, se utilizó una base de datos llamado **Ames Housing Dataset**, obtenido de la plataforma Kaggle. Dicho dataset, muestra un conjunto de 81 columnas, de las cuales 80 son características pertenecientes a 1460 propiedades ubicadas en Ammes, Estados Unidos.

Dentro del set de datos se encontraban variables categóricas (nominales y ordinales) y numéricas (discretas y continuas), a continuación se presenta una tabla con la clasificación y significado de cada una.

Variable	Descripción	Tipo de Dato	Tipo de Variable
Id	Id único del dato	Entero	Número discreto
MSSubClass	Tipo de vivienda	Entero	Categoría nominal
MISoning	Clasificación de la zona	Objeto	Categoría nominal
LotFrontage	Distancia en pies que conecta a la propiedad con la calle	Flotante	Número continuo
LotArea	Área del lote en pies cuadrados	Entero	Número discreto
Street	Calle donde está ubicada la vivienda	Objeto	Categoría nominal
Alley	Tipo de callejón de entrada	Objeto	Categoría nominal
LotShape	Forma general de la propiedad	Objeto	Categoría nominal
LandContour	Ilanura de la propiedad	Objeto	Categoría nominal
Utilities	Utilidades disponibles	Objeto	Categoría nominal
LotConfig	Configuración de lote	Objeto	Categoría nominal
LandSlope	Pendiente	str	Categoría nominal
Neighborhood	Vicindario	Objeto	Categoría nominal
Condition1	Proximidad a calle principal 1	Objeto	Categoría nominal
Condition2	Proximidad a calle principal 2	Objeto	Categoría nominal
BldgType	Tipo de construcción	Objeto	Categoría nominal
HouseStyle	Tipo de vivienda	object	Categoría nominal
OverallQual	Calidad general del material y del acabado	int64	Número Ordinal
OverallCond	Calificación de condición general	int64	Número Ordinal
YearBuilt	Fecha de construcción original	int64	Número Discreta
YearRemodAdd	Año de remodelación	int64	Número Discreta
RoofStyle	Tipo de techo	object	Categoría Nominal
RoofMatl	Material del techo	object	Categoría Nominal
Exterior1st	Revestimiento (material) del exterior	object	Categoría Nominal
Exterior2nd	Revestimiento (material) del exterior, si es más de un material	object	Categoría Nominal
MassnType	Tipo de chapa de mampostería	object	Categoría Nominal
MassnArea	Área de revestimiento de mampostería en pies cuadrados	float64	Número Continuo
ExterQual	Calidad de materiales exteriores	object	Categoría Ordinal

ExterCond	Estado actual del material en el exterior	object	cat ↑ ↓ o o
Foundation	Tipo de cimiento	object	Categoría Nominal
BsmntQual	Alura del sótano	object	Categoría Ordinal
BsmntCond	Estado general del sótano	object	Categoría nominal
BsmntExposure	Muros de sótano a nivel de salida o jardín	object	Categoría Ordinal
Heating	Calefacción	Objeto	Categoría nominal
HeatingQC	Calidad de Calefacción	Objeto	Categoría ordinal
CentralAir	Aire Acondicionado Central	Objeto	Categoría nominal
Electrical	Sistema Eléctrico	Objeto	Categoría nominal
1stFlrSF	Área del Primer Piso	Entero	Número discreto
2ndFlrSF	Área del Segundo Piso	Entero	Número discreto
LowQualFinSF	Área de Acabado de Baja Calidad	Entero	Número discreto
GLWAra	Área de Espacio Habitacional (sobre nivel del suelo)	Entero	Número discreto
BsmntFullBath	Baños Completos en Sótano	Entero	Número discreto
BsmntHalfBath	Medios Baños en Sótano	Entero	Número discreto
FullBath	Baños Completos	Entero	Número discreto
HalfBath	Medios baños	int	Número discreto
BedroomAbvGr	numero de cuartos	int	Número discreto
KitchenAbvGr	numero cocinas	int	Número discreto
KitchenQual	Calidad cocina	object	Categoría ordinal
TotalBsmtAbvGr	Total cuartos	int	Número discreto
Function1	Home functionality	object	Categoría discreta
Fireplaces	Cantidad fireplaces	int	Número discreto
FireplaceQu	calidad fireplace	object	Categoría nominal
GarageType	Ubicación del garage	object	Categoría nominal
GarageYrBl	Años que lleva construido	int	Número continuo
GarageFinish	Acabados interior	object	Categoría nominal
GarageCars	Tamaño garage por cantidad coches	int	Número discreto

GarageFinish	Acabados interior	object	cat ↑ ↓ o o
GarageCars	Tamaño garage por cantidad coches	int	Número discreto
WoodDeckSF	Área de la cubierta	int64	Número discreto
OpenPorchSF	Área del porche	int64	Número discreto
EnclosedPorch	Área de porche cerrado	int64	Número discreto
ScreenPorch	Área de porche de tres estaciones	int64	Número discreto
PoolArea	Área de alberca	int64	Número discreto
PoolQC	Calidad de alberca	object	Categoría ordinal
Fence	Calidad de barda	object	Categoría ordinal
MiscFeature	Características misceláneas	object	Categoría nominal
MiscVal	Valor de las características misceláneas	int64	Número discreto
MoSold	Mes de venta	int64	Categoría nominal
YrSold	Año de venta	int64	Categoría nominal
SaleType	Tipo de venta	object	Categoría nominal
SaleCondition	Condición de venta	object	Categoría nominal
SalePrice	Precio de venta	int64	Número discreto

Figura 1. Fragmento de código de descripción y clasificación de cada variable.

II. PREPROCESAMIENTO

La construcción de cualquier modelo requiere de un proceso de limpieza y análisis previo; preprocesamiento. Este proceso depende fundamentalmente de la naturaleza de los datos con los que se esté trabajando, por lo que cada proceso es diferente y requiere de una explicación y justificación de las medidas tomadas en cada paso.

El preprocesamiento de nuestros datos comienza con un análisis de las características principales de estos. En concreto, poder diferencias entre datos numéricos y datos categóricos, los cuales a su vez estarán agrupados en sus respectivas sub-clasificaciones. Una vez agrupados, se puede realizar diferentes observaciones a estos; de los datos numéricos obtuvimos las medidas de tendencia, mientras que para los categóricos obtuvimos las modas. Para ambos tipos de datos calculamos las distribuciones de frecuencia, lo que cuál es fundamental para tomar decisiones en invitaciones, y para entender la cardinalidad de los mismos.

Posteriormente, se efectuó un trabajo con los datos nulos. Para ese punto es necesario entender que no porque haya apariciones de datos nulos se les puede dar tratamiento igual a las columnas, y es necesario tomar en cuenta el contexto de los datos (lo que miden) para definir qué tratamiento recibirán. Primero, prescindimos de las columnas que presentaban más de 500 datos nulos, lo que sería

más de un tercio de las apariciones. Luego, tiramos las variables que estaban altamente correlacionadas, donde nuestro límite de decisión fue una correlación mayor a 0.80. Este valor es arbitrario, pero es un límite que en convención indica correlaciones considerables. Después, tiramos variables con baja cardinalidad, esto es, variables que tenían muchos de sus registros acomodados en una misma medida. También tiramos variables que podrían propiciar *data leaking*, como lo pueden ser variables del detalle de la compra, las cuales se registran después del precio, el cual es nuestro *target*. Las variables borradas en este poso fueron: 'Street', 'LandContour', 'Utilities', 'LandSlope', 'Condition1', 'Condition2', 'BsmtCond', 'RoofMatl', 'Heating', 'CentralAir', 'Electrical', 'Functional', 'GarageQual', 'BsmtFinType2', 'GarageCond', 'PavedDrive', 'SaleType', 'SaleCondition', 'BldgType', 'ExterCond', 'BsmtFullBath'.

En este punto, tomamos medidas de ingeniería de datos, donde extrajimos información en la columna de *YearRemodAdd* para establecer si una vivienda había sido remodelada sin tener información redundante (año de construcción). En este punto también codificamos las variables ordinales. Esto es importante porque el siguiente paso corresponde a la imputación de valores nulos restantes, donde se usó la moda. Sin embargo, es necesario mencionar que no se imputaron modas en todos los valores nulos, pues muchos de estos correspondían a un valor de escala en el contexto de los valores ordinales y fue necesario codificarlos según ese orden.

El último paso considerable de este preprocesamiento fue una técnica llamada *mean encoder*. La cual agrupa las apariciones únicas por columnas y calcula la media de esas apariciones según su valor de la columna objetivo, en este caso *SalePrice*. Esta columna extrae información del *target*, por lo que debe usarse con cautela. Sin embargo, puede ser apropiada en casos de regresión, donde queremos capturar la relación directa entre la variable y un valor continuo. Finalmente, antes de usar los datos para el entrenamiento, se estandarizaron los datos con la técnica *Min-Max*, de manera que el modelo no se viera sesgado por los datos que naturalmente tiene escalas más grandes.

III. GENERACIÓN DEL MODELO

Una vez limpios los datos, se implementó el modelo de machine learning que permitiría obtener las predicciones más acertadas. Para el entrenamiento local del modelo y la evaluación con *k-fold cross validation* se dividió el conjunto de datos en 70 % de entrenamiento y 30 % de prueba.

Se decidió utilizar un perceptrón multicapa (MLP Regressor) el cual es una red neuronal

artificial que utiliza *backpropagation* para ajustar los pesos entre las neuronas para mejorar la accuracy del modelo.

Dicho modelo es útil tanto para regresiones cuya salida es continua como discreta. Así mismo, su arquitectura consiste en una capa de entrada, capas ocultas y finalmente una capa de salida. En nuestro caso, decidimos implementar 3 capas ocultas, todas de 34 neuronas. Finalmente, como se quiere predecir un solo resultado, la capa de salida es de una sola neurona.

Para la función de activación, se utilizó la función *ReLU*, la cual maneja la no linealidad. Además, se implementó '*LBFGS*' como *solver* para ajustar los pesos. Dicha elección proviene de que es un *solver* que se desenvuelve bastante bien antes muchos parámetros (nosotros tenemos 80 más las columnas de *dummies*) y es usado en regresiones.

Por otra parte, para evitar overfitting en nuestro modelo se utilizó como técnica de regularización '*L2*', con un *alpha* de 1.3 obtenido a través de diversas pruebas con diversos valores.

```

('mlp', MLPRegressor(hidden_layer_sizes=(34, 34, 34),
                      alpha=1.3,
                      activation='relu', solver='lbfgs',
                      random_state=10))

```

Figura 2. Fragmento de código de MLP Regressor

Finalmente, para evaluar el rendimiento de nuestro modelo se utiliza el error cuadrático medio (MSE), la raíz del error cuadrático medio (RMSE) y el R^2 ajustado.

```

R2 score: 0.8722619496266457
MSE: 749002488.5252677
RMSE: 27367.90983113741

```

Figura 3. Resultados del modelo.

Se puede apreciar que el modelo explica con 86% de certeza los datos, por lo cual tiene una confiabilidad bastante grande. Así mismo, su error, comparado con los precios de las casas, es bastante bueno.

Por último, para poder implementar el modelo en nuestro objetivo final, el cual era una aplicación, se debieron de ajustar la cantidad de parámetros a considerar dentro del modelo.

Para el usuario no sería útil ni atractivo el encontrar un cuestionario con las 34 variables finales del modelo. Así mismo, no todas las variables eran fáciles de obtener la información para el usuario, por ejemplo la cantidad de pies que comprendía el recubrimiento de la casa. Por otro lado, había otras variables que eran subjetivas o difíciles de responder, por ejemplo si se esperaba que el usuario clasifique la calidad de su casa, de sus acabados o de un área específica.

Por ende, se optó por, a partir de los resultados de los coeficientes arrojados por el modelo, reducir

aún más su cantidad.

Se optó por dejar 15 variables en el cuestionario y, para no afectar el modelo, las 19 variables más del modelo se mantuvieron con valores estáticos.

Dichas 15 variables debían cumplir con ser de las más significativas para el modelo y de poder ser amigables con el usuario.

Descripción del producto final

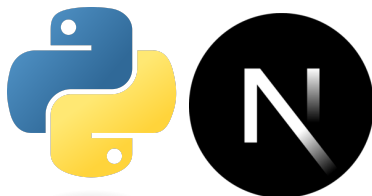
I. ¿CÓMO RESOLVEMOS LA PROBLEMÁTICA?

Describiendo a detalle el modelo previamente presentado, podemos afirmar que este resuelve con creces la problemática inicialmente planteada. En este sentido, ¿cómo transformamos dicha solución técnica en un producto usable y económicamente viable? Así es como concebimos NeighborHub.



Comenzamos definiendo NeighborHub como una aplicación web adaptable, puesto que esto nos concedería presencia en una amplia gama de dispositivos; por ende, una gran cantidad de usuarios.

Una vez establecido el tipo de producto a desarrollar, se eligió Python Flask y Next.js como tecnologías para Backend y Frontend respectivamente.



Así pues, dichas herramientas nos permitieron generar una aplicación profesional, amigable e intuitiva, cuyo propósito principal es permitir al cliente conocer el precio estimado de una propiedad a comprar/vender de la manera más precisa posible.



Conseguimos cumplir el objetivo principal de la plataforma mediante un formulario que recibe las 15 variables establecidas en la descripción del modelo utilizado.

II. MODELO DE NEGOCIOS

Como primer paso de nuestro escalamiento continuo, la plataforma será gratuita para todo tipo de usuarios que deseen predecir el precio de una propiedad a vender o comprar. En este sentido, se tiene prevista la publicidad como medio para generar fondos y así continuar con el desarrollo de la aplicación. Así pues, se estará trabajando exhaustivamente en la creación y optimización de modelos de ML más eficientes y exactos para implementarse.

En este sentido, una vez se cuente con más modelos lo suficientemente robustos, se le ofrecerá al cliente la posibilidad de aportar una pequeña

contribución monetaria a cambio de requests que ofrecerán mucha más precisión.

III. NORMATIVIDAD DE LA SOLUCIÓN IMPLEMENTADA

NeighborHub se compromete a cumplir con todas las leyes y regulaciones en la industria inmobiliaria y tecnológica. Esto incluye:

1. Protección de Datos Personales: Garantizan la privacidad y seguridad de los datos de los usuarios, cumpliendo con leyes como el RGPD. Proporcionan controles para la gestión transparente de datos personales.
2. Transparencia en la Recopilación de Datos: Informan a los usuarios sobre qué datos se recopilan y cómo se utilizan mediante políticas de privacidad claras.
3. No Discriminación: La plataforma es inclusiva y no discrimina a los usuarios por características personales.
4. Exactitud de los resultados: Se esfuerzan por proporcionar estimaciones de precios precisas, minimizando sesgos y utilizando datos confiables.
5. Publicidad Responsable: Compromiso de mostrar anuncios éticos y no invasivos que no afecten la experiencia del usuario.

Cumplimiento de Leyes, Normas y Principios Éticos:

NeighborHub se compromete a cumplir con todas las leyes y regulaciones aplicables en la industria inmobiliaria y en el ámbito de la tecnología de la información. Para garantizar la ética y la transparencia en nuestro servicio, hemos establecido las siguientes medidas:

Normatividad Correspondiente:

A lo largo del proceso de desarrollo, así como el tiempo de operación de la plataforma, hacemos énfasis en la normatividad específica relacionada con la industria inmobiliaria y la tecnología en cada región donde operamos. Esto incluye la correspondiente investigación de regulaciones inmobiliarias locales, estándares de seguridad de datos y cualquier otro requisito legal relevante. Además, proporcionamos información sobre nuestras prácticas que cumplen con estas regulaciones para garantizar la transparencia y la confianza de nuestros usuarios. Estamos comprometidos a mantenernos actualizados y ajustar nuestras operaciones según sea necesario para

cumplir con los cambios en la normativa y los principios éticos de la industria y/o región donde operamos.

Conclusiones.

A partir de lo generado en el transcurso de este proyecto, nos es posible llegar a múltiples conclusiones. Tanto relacionadas con el análisis de datos y el uso de herramientas de inteligencia artificial para su estudio, así como relacionadas con el producto desarrollado a partir de la generación del modelo.

Primeramente, la importancia de la analítica de datos en el mercado inmobiliario. Con lo visto y generado hasta el presente día nos fue posible ver el papel fundamental que el análisis de datos tiene en el mercado inmobiliario; al proporcionar información valiosa para vendedores y compradores. Esta información permite precios justos y eficientes, lo que beneficia a ambas partes y contribuye a una mejora en la experiencia de compra.

Así pues, las agencias inmobiliarias que adoptan la analítica de datos tienen una mayor ventaja competitiva sobre aquellas que no. Es resaltable el enfoque tomado sobre el proyecto para proporcionar datos concretos y predicciones precisas.

Nuestra elección de un perceptrón multicapa como modelo de machine learning es un enfoque efectivo para predecir valores en el mercado inmobiliario. La alta tasa de certeza que tuvo demuestra que la arquitectura de un modelo con capas ocultas y una capa sólida; junto con la función de activación ReLu y el solver LBFGS, en su conjunto, resultan útiles en este contexto en concreto.

Por último, nos parece pertinente señalar que en última instancia la analítica de datos en cualquier industria es un proceso en constante evolución. La mejora continua, tanto en la elección de datos como en la adaptación del modelo a las necesidades cambiantes de los usuarios, es esencial para mantener una ventaja competitiva en este sector; por lo cual, creemos que el trabajo realizado hasta este punto son solo pasos iniciales para un producto en constante crecimiento y mejora.